# Pizza Sales Analysis

**An End-to-End Azure Data Engineering Project**

---

**Objective:**

To analyze pizza sales data by creating a modern, cloud-based data engineering pipeline, leveraging Azure services and Power BI for actionable insights.

---

**Key Features:**

1. **Data Migration**: Seamless transfer of pizza sales data from SQL Server to Azure Blob Storage using Azure Data Factory.

2. **Data Processing:** Advanced transformations and aggregations on Azure Databricks using PySpark and Spark SQL.

3. **Data Visualization:** Insightful dashboards in Power BI showcasing KPIs such as total sales, pizza popularity, and revenue trends.

---

**Technologies Used:**

- **SQL Server:** Source database.

- **Azure Data Factory:** Data pipeline orchestration.

- **Azure Blob Storage:** Scalable cloud storage.

- **Azure Databricks:** Data transformation and aggregation.

- **Power BI:** Visualization and reporting.

---

**Prepared By:**

## Alkesh Lajurkar

**LinkedIn:** linkedin.com/in/alkeshlajurkar
**GitHub:** github.com/alkeshlajurkar
**Email:** alkeshlajurkar@gmail.com
**Contact:** +91 8390956095

---

# 1. Introduction

This project demonstrates how to design and implement an end-to-end data engineering pipeline. The pipeline includes migrating data from an on-premises SQL Server to Azure, processing it with Azure Databricks, and visualizing insights in Power BI. This comprehensive guide walks through the steps to implement the solution while ensuring scalability, automation, and actionable insights.

**Skills Demonstrated**

- SQL Server

- Azure Data Factory

- Azure Databricks (PySpark & Spark SQL)

- Power BI

- End-to-End Data Engineering Pipeline

---

# 2. Business Requirement Overview

**Objective**

The goal is to build a pipeline that efficiently migrates and processes data from on-prem SQL Server to Azure for business reporting. The final output is a Power BI dashboard presenting KPIs and trends derived from the data.

**Tasks Overview**

1. Import raw sales data from CSV to SQL Server.

2. Migrate data from SQL Server to Azure Blob Storage using Azure Data Factory.

3. Process and transform the data in Azure Databricks.

4. Visualize the processed data in Power BI.

**Tools Used**

- **SQL Server:** On-prem database for raw data.

- **Azure Storage Account:** Cloud storage to hold migrated data.

- **Azure Data Factory (ADF):** Tool for orchestrating data movement.

- **Azure Databricks:** For data processing and KPI generation.

- **Power BI:** Visualization tool for creating dashboards.

---

# 3. Step-by-Step Implementation

## Step 1: Import CSV Data to SQL Server

**Objective:** Load raw sales data into a SQL Server database for processing.

1. **Create the Database:** Use SQL Server Management Studio (SSMS) to create a database.

2. **Create the Table:** Define a table named PizzaSales with columns to store sales data.

3. **Import Data from CSV:**

   o Use SSMS's *Import Flat File Wizard* to upload data from a CSV file.

   o Map columns correctly and execute the process.

4. **Verify Data:** Run a query to ensure the data was imported successfully.



## Step 2: Set Up Azure Resources

**Objective:** Prepare Azure services for cloud-based data storage and processing.

1. **Create a Storage Account:**

   o Go to the Azure portal.

   o Create a storage account under your resource group.

   o Configure settings for secure data storage.

2. **Create a Blob Container:**

- o Inside the storage account, create a new container named salesdata.
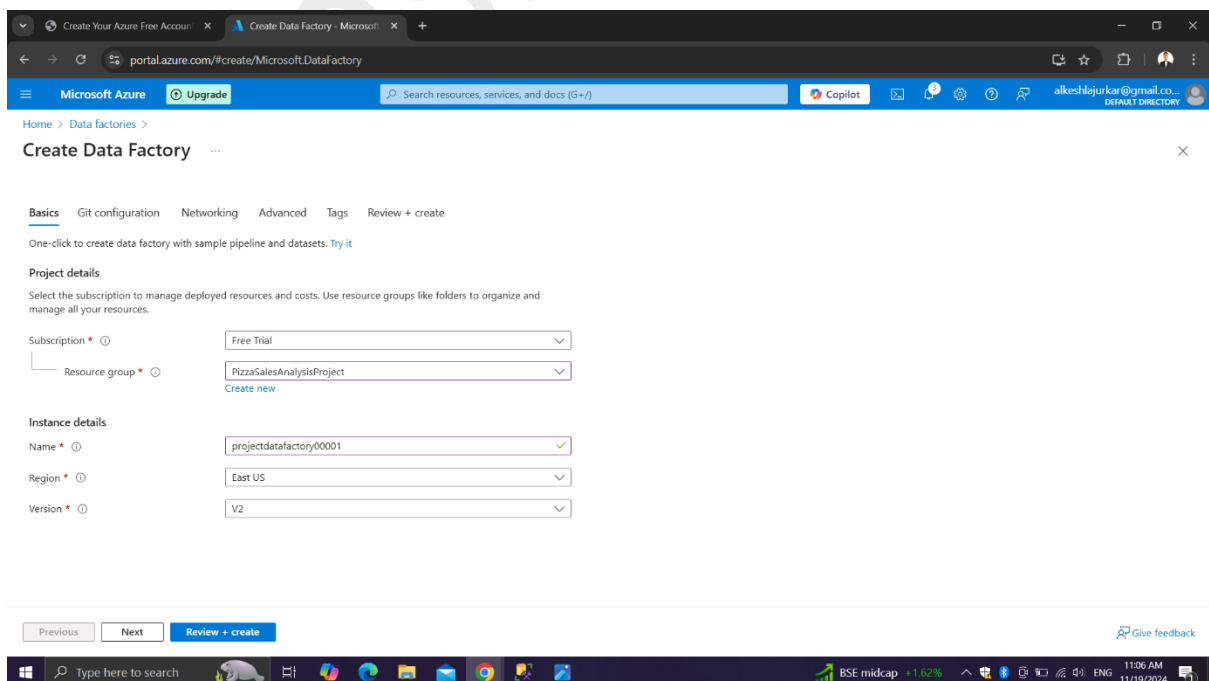
- o This will store the raw files after migration.

## Step 3: Migrate Data Using Azure Data Factory

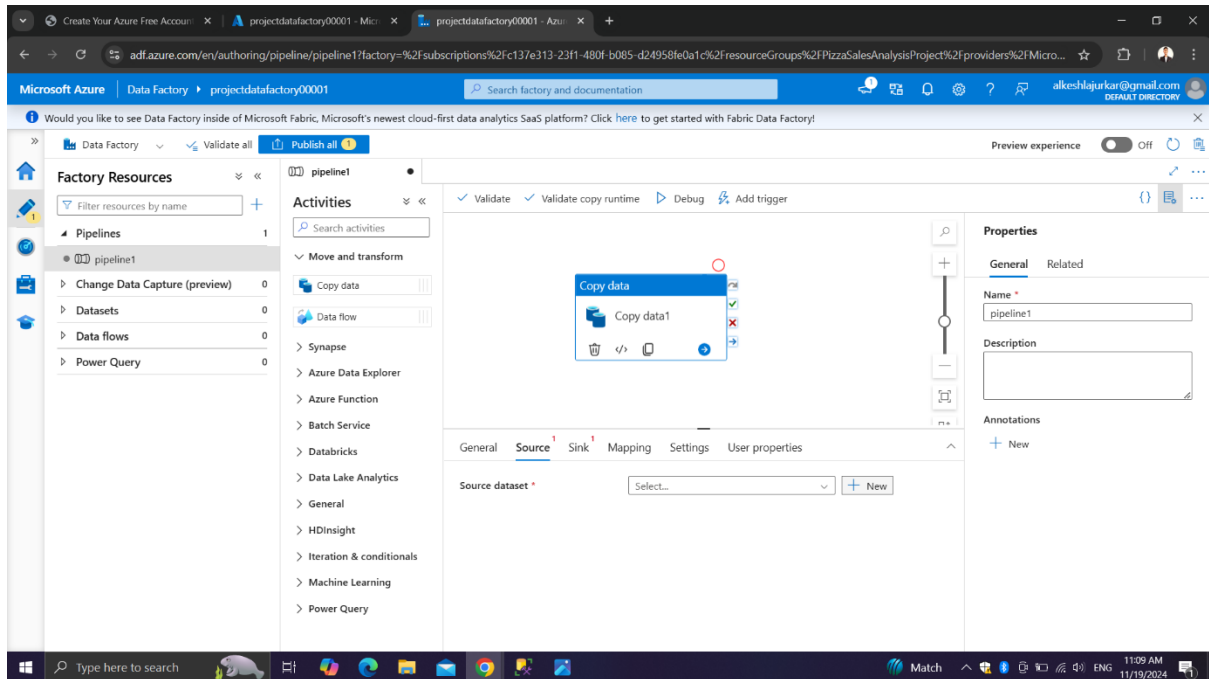**Objective:** Automate the transfer of data from SQL Server to Azure Blob Storage.

1. **Create an Azure Data Factory Instance:**

   o   In the Azure portal, create a new Azure Data Factory instance.
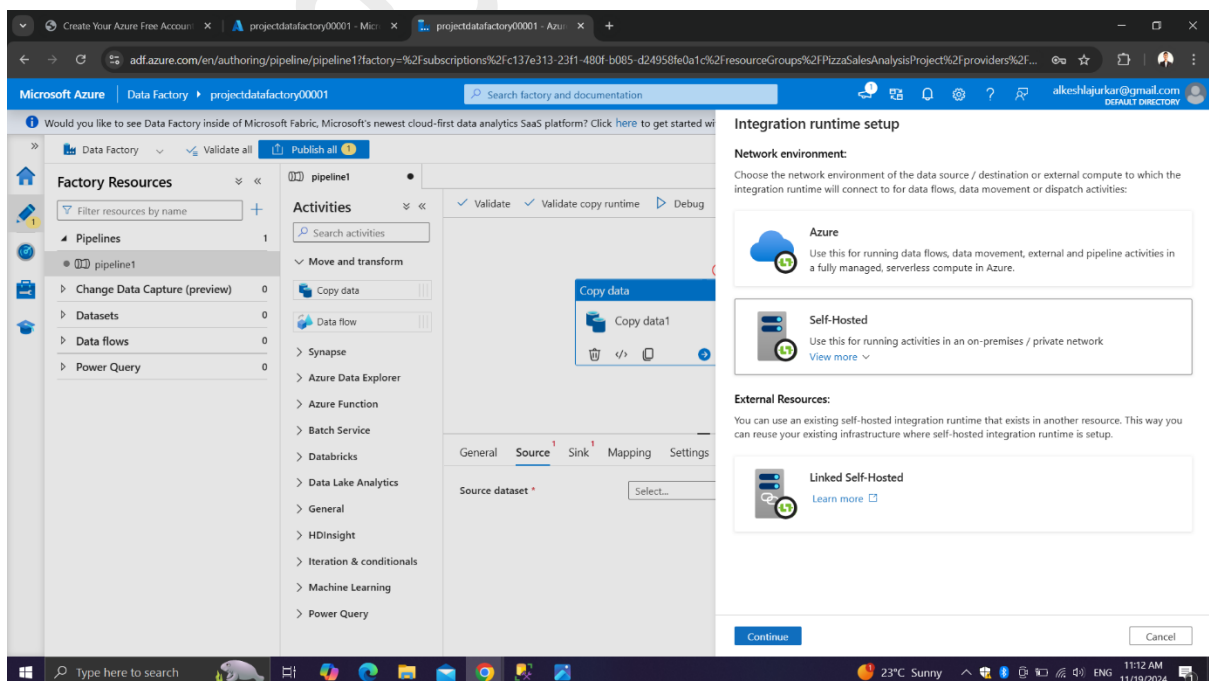
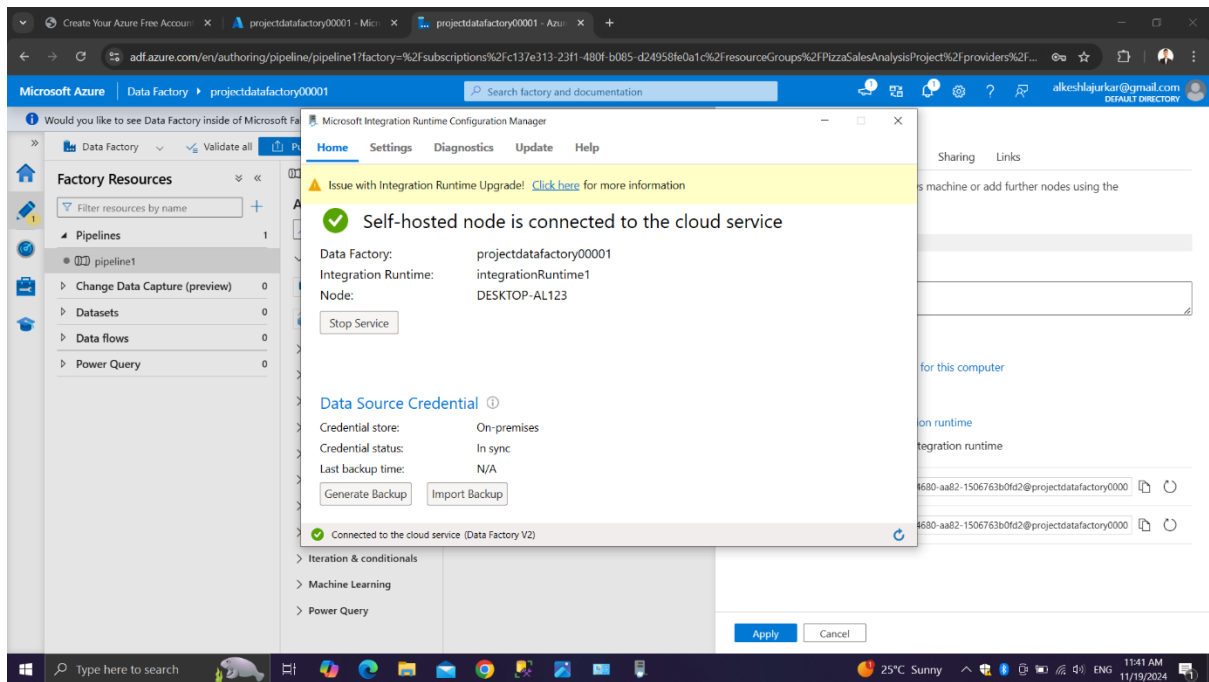   o   Name it appropriately and configure the region.

## 2. Build the Pipeline:

- o Use the *Copy Data Activity* to create a pipeline that transfers data from SQL Server to Blob Storage.

- o Map source (SQL Server) and destination (Blob Storage).



## 3. Configure Integration Runtime:

- o Set up a self-hosted integration runtime to securely connect on-prem SQL Server to Azure.
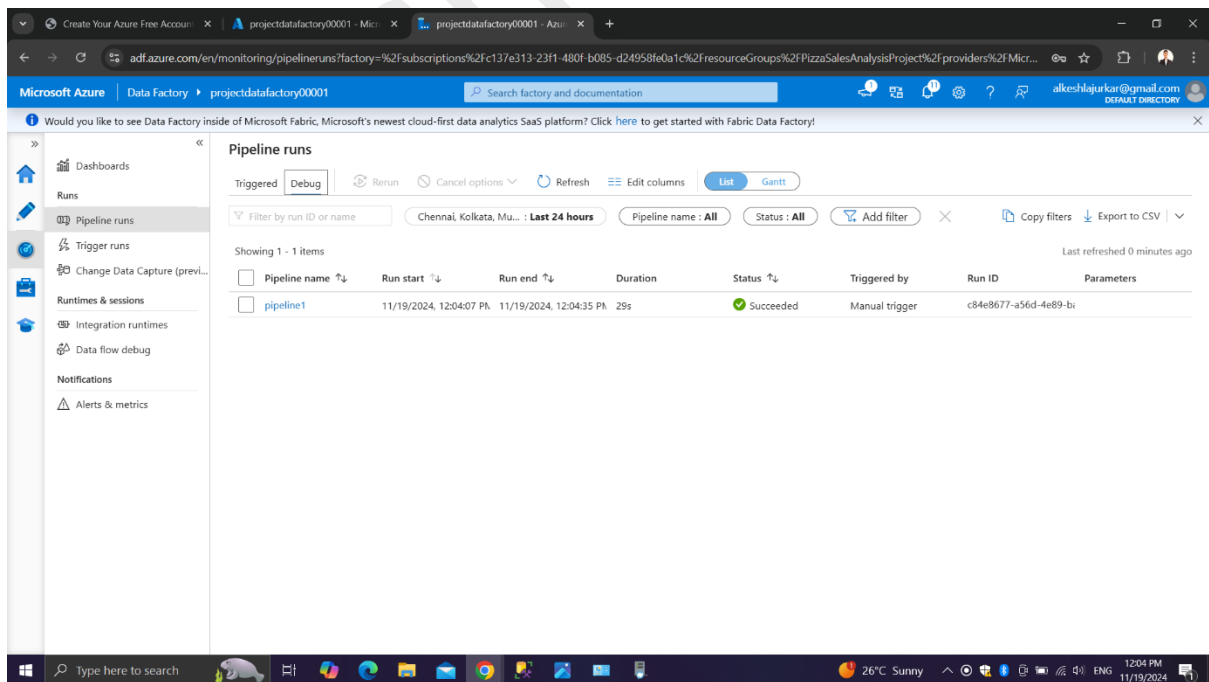
## 4. Set Up Linked Services:

- **SQL Server Linked Service:** Connect to your on-prem SQL Server.

- **Blob Storage Linked Service:** Connect to your Azure Blob Storage.

### 5. Run and Verify the Pipeline:

- o Execute the pipeline and verify that the raw file has been successfully transferred to the salesdata container in Blob Storage.
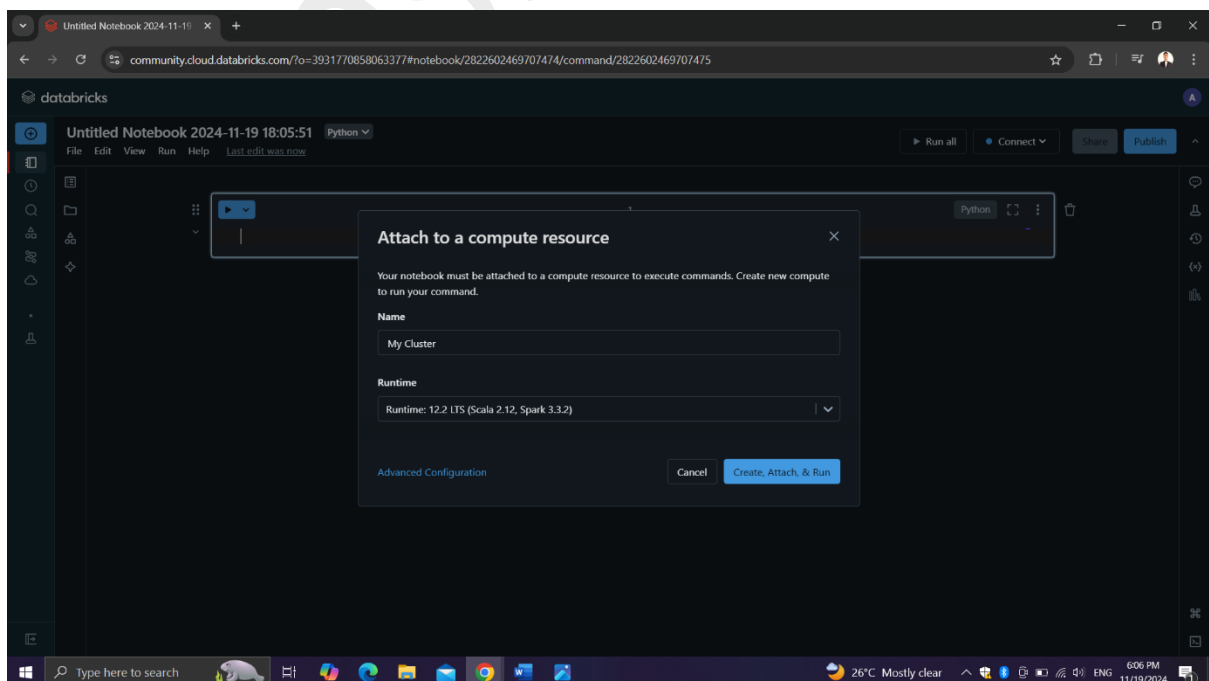
## Step 4: Data Processing in Azure Databricks

**Objective:** Process raw data and generate KPIs using Spark-based workflows.

1. **Set Up Azure Databricks Workspace:**

   o In the Azure portal, create a new Databricks workspace.

   o Launch the workspace and create a new cluster.

2. **Mount Blob Storage in Databricks:**

   o Use the following PySpark code to mount the salesdata container in Databricks.



3. **Read Data into a PySpark DataFrame:**

4. **Perform Data Transformations:**

   o Calculate KPIs such as total sales and total quantity.

5. **Save the Transformed Data:**

   o Write the processed data back to Blob Storage for Power BI consumption.

---

**Step 5: Visualize Data in Power BI**

**Objective:** Create interactive dashboards to present insights.

1. **Connect Power BI to Azure Databricks:**

   o Install the Databricks ODBC driver.

   o Use the driver to connect Power BI to Databricks.

2. **Build Dashboards:**

   o Create visuals such as:

     ▪ Bar charts for total sales by pizza type.

     ▪ Pie charts showing quantity sold distribution.

     ▪ Line graphs for sales trends over time.

   o Use slicers for filtering data by date or product.

3. **Publish Reports:**

   o Publish the Power BI report to the Power BI service for collaboration.

## 4. Outcome

**Project Deliverables**

1. **SQL Server to Azure Data Migration:** Successfully moved on-prem data to Blob Storage using Azure Data Factory.

2. **Data Transformation in Azure Databricks:** Processed raw data into meaningful insights.

3. **Power BI Dashboards:** Developed interactive dashboards showcasing KPIs like total sales and quantity sold.

**Key Benefits**

- **Automation:** Eliminated manual intervention by setting up a scheduled pipeline.

- **Scalability:** Leveraged Azure's robust infrastructure for future growth.

- **Actionable Insights:** Delivered real-time business metrics for informed decision-making.

## 5. Conclusion

This project showcases a robust end-to-end Azure data engineering pipeline. By combining the capabilities of SQL Server, Azure Data Factory, Azure Databricks, and Power BI, it provides a scalable solution for data migration, processing, and visualization. These skills are essential for modern data engineers tackling real-world challenges.